

Lagrangian Relaxation Applied to Sparse Global Network Alignment

Mohammed El-Kebir^{1,2,3}, Jaap Heringa^{2,3,4}, and Gunnar W. Klau^{1,3}

¹ Centrum Wiskunde & Informatica, Life Sciences Group, Science Park 123, 1098 XG Amsterdam, the Netherlands, {m.el-kebir,gunnar.klau}@cwi.nl

² Centre for Integrative Bioinformatics VU (IBIVU), VU University Amsterdam, De Boelelaan 1081A, 1081 HV Amsterdam, the Netherlands, heringa@few.vu.nl

³ Netherlands Institute for Systems Biology, Amsterdam, the Netherlands

⁴ Netherlands Bioinformatics Centre, Nijmegen, the Netherlands

Abstract. Data on molecular interactions is increasing at a tremendous pace, while the development of solid methods for analyzing this network data is lagging behind. This holds in particular for the field of comparative network analysis, where one wants to identify commonalities between biological networks. Since biological functionality primarily operates at the network level, there is a clear need for topology-aware comparison methods. In this paper we present a method for global network alignment that is fast and robust, and can flexibly deal with various scoring schemes taking both node-to-node correspondences as well as network topologies into account. It is based on an integer linear programming formulation, generalizing the well-studied quadratic assignment problem. We obtain strong upper and lower bounds for the problem by improving a Lagrangian relaxation approach and introduce the software tool NATALIE 2.0, a publicly available implementation of our method. In an extensive computational study on protein interaction networks for six different species, we find that our new method outperforms alternative state-of-the-art methods with respect to quality and running time.

1 Introduction

In the last decade, data on molecular interactions has increased at a tremendous pace. For instance, the STRING database [24], which contains protein protein interaction (PPI) data, grew from 261,033 proteins in 89 organisms in 2003 to 5,214,234 proteins in 1,133 organisms in May 2011, more than doubling the number of proteins in the database every two years. The same trends can be observed for other types of biological networks, including metabolic, gene-regulatory, signal transduction and metagenomic networks, where the latter can incorporate the excretion and uptake of organic compounds through, for example, a microbial community [12, 21]. In addition to the plethora of experimentally derived network data for many species, also the structure and behavior of molecular networks have become intensively studied over the last few years [2], leading to the observation of many conserved features at the network level. However, the

development of solid methods for analyzing network data is lagging behind, particularly in the field of comparative network analysis. Here, one wants to identify commonalities between biological networks from different strains or species, or derived from different conditions. Based on the assumption that evolutionary conservation implies functional significance, comparative approaches may help (i) improve the accuracy of data, (ii) generate, investigate, and validate hypotheses, and (iii) transfer functional annotations. Until recently, the most common way of comparing two networks has been to solely consider node-to-node correspondences, for example by finding homologous relationships between nodes (e.g. proteins in PPI networks) of either network, while the topology of the two networks has not been taken into account. Since biological functionality primarily operates at the network level, there is a clear need for topology-aware comparison methods. In this paper we present a network alignment method that is fast and robust, and can flexibly deal with various scoring schemes taking both node-to-node correspondences as well as network topologies into account.

Previous work. Network alignment establishes node correspondences based on both node-to-node similarities and conserved topological information. Similar to sequence alignment, *local* network alignment aims at identifying one or more shared subnetworks, whereas *global* network alignment addresses the overall comparison of the complete input networks.

Over the last years a number of methods have been proposed for both global and local network alignment, for example PATHBLAST [14], NETWORKBLAST [22], MAWISH [16], GRAEMLIN [8], ISORANK [23], GRAAL [17], and SUBMAP [4]. PATHBLAST heuristically computes high-scoring similar paths in two PPI networks. Detecting protein complexes has been addressed with NETWORKBLAST by Sharan et al. [22], where the authors introduce a probabilistic model and propose a heuristic greedy approach to search for shared complexes. Koyutürk et al. [16] use a more elaborate scoring scheme based on an evolutionary model to compute local pairwise alignments of PPI networks. The ISORANK algorithm by Singh et al. [23] approaches the global alignment problem by preferably matching nodes which have a similar neighborhood, which is elegantly solved as an eigenvalue problem. Kuchaiev et al. [17] take a similar approach. Their method GRAAL matches nodes that share a similar distribution of so-called graphlets, which are small connected non-isomorphic induced subgraphs.

In this paper we focus on pairwise global network alignment, where an alignment is scored by summing up individual scores of aligned node and interaction pairs. Among the above mentioned methods, ISORANK and GRAAL use a scoring model that can be expressed in this manner.

Contribution. We present an algorithm for global network alignment based on an integer linear programming (ILP) formulation, generalizing the well-studied quadratic assignment problem (QAP). We improve upon an existing Lagrangian relaxation approach presented in previous work [15] to obtain strong upper and lower bounds for the problem. We exploit the closeness to QAP and generalize

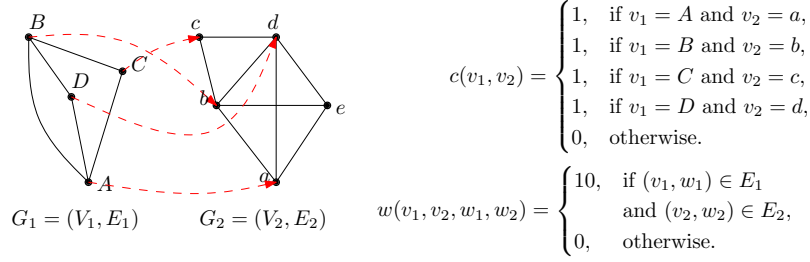


Fig. 1: Example of a network alignment. With the given scoring function, the alignment has a score of $4 + 40 = 44$.

a dual descent method for updating the Lagrangian multipliers to the generalized problem. We have implemented the revised algorithm from scratch as the software tool NATALIE 2.0. In an extensive computational study on protein interaction networks for six different species, we compare NATALIE 2.0 to GRAAL and ISORANK, evaluating the number of conserved edges as well as functional coherence of the modules in terms of GO annotation. We find that NATALIE 2.0 outperforms the alternative methods with respect to quality and running time. Our software tool NATALIE 2.0 as well as all data sets used in this study are publicly available at <http://planet-lisa.net>.

2 Preliminaries

Given two simple graphs $G_1 = (V_1, E_1)$ and $G_2 = (V_2, E_2)$, an *alignment* $a : V_1 \rightarrow V_2$ is a *partial injective function* from V_1 to V_2 . As such we have that an alignment relates every node in V_1 to at most one node in V_2 and that conversely every node in V_2 has at most one counterpart in V_1 . An alignment is assigned a real-valued *score* using an additive scoring function s defined as follows:

$$s(a) = \sum_{v \in V_1} c(v, a(v)) + \sum_{\substack{v, w \in V_1 \\ v < w}} w(v, a(v), w, a(w)) \quad (1)$$

where $c : V_1 \times V_2 \rightarrow \mathbb{R}$ is the score of aligning a pair of nodes in V_1 and V_2 respectively. On the other hand, $w : V_1 \times V_2 \times V_1 \times V_2 \rightarrow \mathbb{R}$ allows for scoring topological similarity. The problem of global pairwise network alignment (GNA) is to find the highest scoring alignment a^* , i.e. $a^* = \arg \max s(a)$. Figure 1 shows an example.

NP-hardness of GNA follows by a simple reduction from the decision problem CLIQUE, which asks whether there is a clique of cardinality at least k in a given simple graph $G = (V, E)$ [13]. The corresponding GNA instance concerns the alignment of the complete graph of k vertices $K_k = (V_k, E_k)$ with G using the scoring function $s(a) = |\{(v, w) \in E_k \mid (a(v), a(w)) \in E\}|$. Since an alignment is injective, there is a clique of cardinality at least k if and only if

the cost of the optimal alignment is $\binom{k}{2}$. The close relationship of GNA with the quadratic assignment problem is more easily observed when formulating GNA as a mathematical program. Throughout the remainder of the text we use dummy variables $i, j \in \{1, \dots, |V_1|\}$ and $k, l \in \{1, \dots, |V_2|\}$ to denote nodes in V_1 and V_2 , respectively. Let C be a $|V_1| \times |V_2|$ matrix such that $c_{ik} = c(i, k)$ and let W be a $(|V_1| \times |V_2|) \times (|V_1| \times |V_2|)$ matrix whose entries w_{ikjl} correspond to interaction scores $w(i, k, j, l)$. Now we can formulate GNA as

$$\max_x \quad \sum_{i,k} c_{ik} x_{ik} + \sum_{\substack{i,j \\ i < j}} \sum_{\substack{k,l \\ k \neq l}} w_{ikjl} x_{ik} x_{jl} \quad (\text{IQP})$$

$$\text{s.t.} \quad \sum_l x_{jl} \leq 1 \quad \forall j \quad (2)$$

$$\sum_j x_{jl} \leq 1 \quad \forall l \quad (3)$$

$$x_{ik} \in \{0, 1\} \quad \forall i, k \quad (4)$$

where the decision variable x_{ik} indicates whether the i -th node in V_1 is aligned with the k -th node in V_2 . The above formulation shares many similarities with Lawler's formulation [19] of the QAP. However, instead of finding an assignment we are interested in finding a matching, which is reflected in constraints (2) and (3) being inequalities rather than equalities. As can be seen in (1) we only consider the upper triangle of W rather than the entire matrix. An analogous way of looking at this, is to consider W to be symmetric. This is usually not the case for QAP instances. In addition, due to the fact that biological input graphs are typically sparse, we have that W is sparse as well. These differences allow us to come up with an effective method of solving the problem as we will see in the following.

3 Method

The relaxation presented here follows the same lines as the one given by Adams and Johnson for the QAP [1]. We start by linearizing (IQP) by introducing binary variables y_{ikjl} defined as $y_{ikjl} := x_{ik} x_{jl}$ and constraints $y_{ikjl} \leq x_{jl}$ and $y_{ikjl} \leq x_{ik}$ for all $i \leq j$ and $k \neq l$. If we assume that all entries in W are positive, we do not need to enforce that $y_{ikjl} \geq x_{ik} + x_{jl} - 1$. In Section 5 we will discuss this assumption. Rather than using the aforementioned constraints, we make use of a stronger set of constraints which we obtain by multiplying constraints (2) and (3) by x_{ik} :

$$\sum_{\substack{l \\ l \neq k}} y_{ikjl} = \sum_{\substack{l \\ l \neq k}} x_{ik} x_{jl} \leq \sum_l x_{ik} x_{jl} \leq x_{ik}, \quad \forall i, j, k, i < j \quad (5)$$

$$\sum_{\substack{j \\ j > i}} y_{ikjl} = \sum_{\substack{j \\ j > i}} x_{ik} x_{jl} \leq \sum_j x_{ik} x_{jl} \leq x_{ik}, \quad \forall i, k, l, k \neq l \quad (6)$$

We proceed by splitting the variable y_{ikjl} (where $i < j$ and $k \neq l$). In other words, we extend the objective function such that the counterpart of y_{ikjl} becomes y_{jlik} . This is accomplished by rewriting the dummy constraint in (6) to $j \neq i$. In addition, we split the weights: $w_{ikjl} = w_{jlik} = (w'_{ikjl}/2)$ where w'_{ikjl} denotes the original weight. Furthermore, we require that the counterparts of the split decision variables assume the same value, which amounts to

$$\max_{x,y} \quad \sum_{i,k} c_{ik} x_{ik} + \sum_{\substack{i,j \\ i < j}} \sum_{\substack{k,l \\ k \neq l}} w_{ikjl} y_{ikjl} + \sum_{\substack{i,j \\ i > j}} \sum_{\substack{k,l \\ k \neq l}} w_{ikjl} y_{ikjl} \quad (\text{ILP})$$

$$\text{s.t.} \quad \sum_l x_{jl} \leq 1 \quad \forall j \quad (7)$$

$$\sum_j x_{jl} \leq 1 \quad \forall l \quad (8)$$

$$\sum_{\substack{l \\ l \neq k}} y_{ikjl} \leq x_{ik} \quad \forall i, j, k, i \neq j \quad (9)$$

$$\sum_{\substack{j \\ j \neq i}} y_{ikjl} \leq x_{ik} \quad \forall i, k, l, k \neq l \quad (10)$$

$$y_{ikjl} = y_{jlik} \quad \forall i, j, k, l, i < j, k \neq l \quad (11)$$

$$y_{ikjl} \in \{0, 1\} \quad \forall i, j, k, l, i \neq j, k \neq l \quad (12)$$

$$x_{ik} \in \{0, 1\} \quad \forall i, k \quad (13)$$

We can solve the continuous relaxation of (ILP) via its Lagrangian dual by dualizing the linking constraints (11) with multiplier λ :

$$\min_{\lambda} \quad Z_{\text{LD}}(\lambda) , \quad (\text{LD})$$

where $Z_{\text{LD}}(\lambda)$ equals

$$\begin{aligned} \max_{x,y} \quad & \sum_{i,k} c_{ik} x_{ik} + \sum_{\substack{i,j \\ i < j}} \sum_{\substack{k,l \\ k \neq l}} (w_{ikjl} + \lambda_{ikjl}) y_{ikjl} + \sum_{\substack{i,j \\ i > j}} \sum_{\substack{k,l \\ k \neq l}} (w_{ikjl} - \lambda_{jlik}) y_{ikjl} \\ \text{s.t.} \quad & (7), (8), (9), (10), (12) \text{ and } (13) \end{aligned}$$

Now that the linking constraints have been dualized, one can observe that the remaining constraints decompose the variables into $|V_1||V_2|$ disjoint groups, where variables across groups are not linked by any constraint, and where each group

contains a variable x_{ik} and variables y_{ikjl} for $j \neq i$ and $l \neq k$. Hence, we have

$$Z_{LD}(\lambda) = \max_x \sum_{i,k} [c_{ik} + v_{ik}(\lambda)] x_{ik} \quad (LD_\lambda)$$

$$\text{s.t.} \quad \sum_l x_{jl} \leq 1 \quad \forall j \quad (14)$$

$$\sum_j x_{jl} \leq 1 \quad \forall l \quad (15)$$

$$x_{ik} \in \{0, 1\} \quad \forall i, k \quad (16)$$

which corresponds to a maximum weight bipartite matching problem on the so-called *alignment graph* $G_m = (V_1 \cup V_2, E_m)$. In the general case G_m is a complete bipartite graph, i.e. $E_m = \{(i, k) \mid i \in V_1, v_2 \in V_2\}$. However, by exploiting biological knowledge one can make G_m more sparse by excluding biologically-unlikely edges (see Section 4). For the global problem, the weight of a matching edge (i, k) is set to $c_{ik} + v_{ik}(\lambda)$, where the latter term is computed as

$$v_{ik}(\lambda) = \max_y \sum_{\substack{j \\ j > i}} \sum_{\substack{l \\ l \neq k}} (w_{ikjl} + \lambda_{ikjl}) y_{ikjl} + \sum_{\substack{j \\ j < i}} \sum_{\substack{l \\ l \neq k}} (w_{ikjl} - \lambda_{jlik}) y_{ikjl} \quad (LD_\lambda^{ik})$$

$$\text{s.t.} \quad \sum_{\substack{l \\ l \neq k}} y_{ikjl} \leq 1 \quad \forall j, j \neq i \quad (17)$$

$$\sum_{\substack{j \\ j \neq i}} y_{ikjl} \leq 1 \quad \forall l, l \neq k \quad (18)$$

$$y_{ikjl} \in \{0, 1\} \quad \forall j, l. \quad (19)$$

Again, this is a maximum weight bipartite matching problem on the same alignment graph but excluding edges incident to either i or k and using different edge weights: the weight of an edge (j, l) is $w_{ikjl} + \lambda_{ikjl}$ if $j > i$, or $w_{ikjl} - \lambda_{jlik}$ if $j < i$. So in order to compute $Z_{LD}(\lambda)$, we need to solve a total number of $|V_1||V_2| + 1$ maximum weight bipartite matching problems, which, using the Hungarian algorithm [18, 20] can be done in $O(n^5)$ time, where $n = \max(|V_1|, |V_2|)$. In case the alignment graph is sparse, i.e. $O(|E_m|) = O(n)$, $Z_{LD}(\lambda)$ can be computed in $O(n^4 \log n)$ time using the successive shortest path variant of the Hungarian algorithm [7]. It is important to note that for any λ , $Z_{LD}(\lambda)$ is an upper bound on the score of an optimal alignment. This is because any alignment a is feasible to $Z_{LD}(\lambda)$ and does not violate the original linking constraints and therefore has an objective value equal to $s(a)$. In particular, the optimal alignment a^* is also feasible to $Z_{LD}(\lambda)$ and hence $a^* \leq Z_{LD}(\lambda)$. Since the two sets of problems resulting from the decomposition both have the integrality property [6], the smallest upper bound we can achieve equals the linear programming (LP) bound of the continuous relaxation of (ILP) [9]. However, computing the smallest upper bound by finding suitable multipliers is much faster than solving the corresponding LP. Given solution (x, y) to $Z_{LD}(\lambda)$, we obtain a lower bound on $s(a^*)$, denoted $Z_{lb}(\lambda)$, by considering the score of the alignment encoded in x .

3.1 Solving Strategies

In this section we will discuss strategies for identifying Lagrangian multipliers λ that yield an as small as possible gap between the upper and lower bound resulting from the solution to $Z_{LD}(\lambda)$.

Subgradient optimization. We start by discussing subgradient optimization, which is originally due to Held and Karp [10]. The idea is to generate a sequence $\lambda^0, \lambda^1, \dots$ of Lagrangian multiplier vectors starting from $\lambda^0 = \mathbf{0}$ as follows:

$$\lambda_{ikjl}^{t+1} = \lambda_{ikjl}^t - \frac{\alpha \cdot (Z_{LD}(\lambda) - Z_{lb}(\lambda))}{\|g(\lambda^t)\|^2} g(\lambda_{ikjl}^t) \quad \forall i, j, k, l, \quad i < j, k \neq l \quad (20)$$

where $g(\lambda_{ikjl}^t)$ corresponds to the subgradient of multiplier λ_{ikjl}^t , i.e. $g(\lambda_{ikjl}^t) = y_{ikjl} - y_{jlik}$, and α is the step size parameter. Initially α is set to 1 and it is halved if neither $Z_{LD}(\lambda)$ nor $Z_{lb}(\lambda)$ have improved for over N consecutive iterations. Conversely, α is doubled if M times in a row there was an improvement in either $Z_{LD}(\lambda)$ or $Z_{lb}(\lambda)$ [5]. In case all subgradients are zero, the optimal solution has been found and the scheme terminates. Note that this is not guaranteed to happen. Therefore we abort the scheme after exceeding a time limit or a pre-specified number of iterations. In addition, we terminate if α has dropped below machine precision. Algorithm 1 gives the pseudo code of this procedure.

Algorithm 1: SUBGRADIENTOPT(λ, M, N)

```

1  $\alpha \leftarrow 1; n \leftarrow N; m \leftarrow M$ 
2  $[LB^*, UB^*] \leftarrow [Z_{lb}(\lambda), Z_{LD}(\lambda)]$ 
3 while  $g(\lambda) \neq 0$  do
4    $\lambda \leftarrow \lambda - \frac{\alpha(Z_{LD}(\lambda) - Z_{lb}(\lambda))}{\|g(\lambda^t)\|^2} g(\lambda^t)$ 
5   if  $[LB^*, UB^*] \setminus [Z_{lb}(\lambda), Z_{LD}(\lambda)] = \emptyset$  then  $n \leftarrow n - 1$ 
6   else
7      $LB^* \leftarrow \max[LB^*, Z_{lb}(\lambda)]$ 
8      $UB^* \leftarrow \min[UB^*, Z_{LD}(\lambda)]$ 
9      $m \leftarrow m - 1$ 
10  if  $n = 0$  then  $\alpha \leftarrow \alpha/2; n \leftarrow N$ 
11  if  $m = 0$  then  $\alpha \leftarrow 2\alpha; m \leftarrow M$ 
12 return  $[LB^*, UB^*]$ 

```

Dual descent. In this section we derive a dual descent method which is an extension of the one presented in [1]. The dual descent method takes as a starting

point the dual of $Z_{LD}(\lambda)$:

$$Z_{LD}(\lambda) = \min_{\alpha, \beta} \sum_i \alpha_i + \sum_k \beta_k \quad (21)$$

$$\text{s.t. } \alpha_i + \beta_k \geq c_{ik} + v_{ik}(\lambda) \quad \forall i, k \quad (22)$$

$$\alpha_i \geq 0 \quad \forall i \quad (23)$$

$$\beta_k \geq 0 \quad \forall k \quad (24)$$

where the dual of $v_{ik}(\lambda)$ is

$$v_{ik}(\lambda) = \min_{\mu, \nu} \sum_{\substack{j \\ j \neq i}} \mu_j^{ik} + \sum_{\substack{l \\ l \neq k}} \nu_l^{ik} \quad (25)$$

$$\text{s.t. } \mu_j^{ik} + \nu_l^{ik} \geq w_{ikjl} + \lambda_{ikjl} \quad \forall j, l, j > i, l \neq k \quad (26)$$

$$\mu_j^{ik} + \nu_l^{ik} \geq w_{ikjl} - \lambda_{jlik} \quad \forall j, l, j < i, l \neq k \quad (27)$$

$$\mu_j^{ik} \geq 0 \quad \forall j \quad (28)$$

$$\nu_l^{ik} \geq 0 \quad \forall l. \quad (29)$$

Suppose that for a given λ^t we have computed dual variables (α, β) solving (21) with objective value $Z_{LD}(\lambda^t)$, as well as dual variables (μ^{ik}, ν^{ik}) yielding values $v_{ik}(\lambda)$ to linear programs (25). The goal now is to find λ^{t+1} such that the resulting bound is better or just as good, i.e. $Z_{LD}(\lambda^{t+1}) \leq Z_{LD}(\lambda^t)$. We prevent the bound from increasing, by ensuring that the dual variables (α, β) remain feasible to (21). This we can achieve by considering the slacks: $\pi_{ik}(\lambda) = \alpha_i + \beta_k - c_{ik} - v_{ik}(\lambda)$. So for (α, β) to remain feasible, we can only allow every $v_{ik}(\lambda^t)$ to increase by as much as $\pi_{ik}(\lambda^t)$. We can achieve such an increase by considering linear programs (25) and their slacks defined as

$$\gamma_{ikjl}(\lambda) = \begin{cases} \mu_j^{ik} + \nu_l^{ik} - w_{ikjl} + \lambda_{ikjl}, & \text{if } j > i, \\ \mu_j^{ik} + \nu_l^{ik} - w_{ikjl} - \lambda_{jlik}, & \text{if } j < i, \end{cases} \quad \forall j, l, j \neq i, l \neq k, \quad (30)$$

and update the multipliers in the following way.

Lemma 1. *The adjustment scheme below yields solutions to linear programs (25) with objective values $v_{ik}(\lambda^{t+1})$ at most $\pi_{ik}(\lambda^t) + v_{ik}(\lambda^t)$ for all i, k .*

$$\begin{aligned} \lambda_{ikjl}^{t+1} = & \lambda_{ikjl}^t + \varphi_{ikjl} \left[\gamma_{ikjl}(\lambda^t) + \tau_{ik} \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t) \right] \\ & - \varphi_{jlik} \left[\gamma_{jlik}(\lambda^t) + \tau_{jl} \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{jl}(\lambda^t) \right] \end{aligned} \quad (31)$$

for all $j, l, i < j, k \neq l$, where $n_1 = |V_1|$, $n_2 = |V_2|$, and $0 \leq \varphi_{ikjl}, \tau_{jl} \leq 1$ are parameters.

Proof. We prove the lemma by showing that for any i, k there exists a feasible solution (μ'^{ik}, ν'^{ik}) to (25) whose objective value $v_{ik}(\lambda^{t+1})$ is at most $\pi_{ik}(\lambda^t) +$

$v_{ik}(\lambda^t)$. Let (μ^{ik}, ν^{ik}) be the solution to (25) given multipliers λ^t . We claim that setting

$$\begin{aligned}\mu_j'^{ik} &= \mu_j^{ik} + \frac{\pi_{ik}(\lambda^t)}{2(n_1 - 1)} & \forall j, j \neq i \\ \nu_l'^{ik} &= \nu_l^{ik} + \frac{\pi_{ik}(\lambda^t)}{2(n_2 - 1)} & \forall l, l \neq k,\end{aligned}$$

results in a feasible solution to (25) given multipliers λ^{t+1} . We start by showing that constraints (26) and (27) are satisfied. From (31) the following bounds on λ^{t+1} follow.

$$\begin{aligned}\lambda_{ikjl}^t - \gamma_{jlik}(\lambda^t) - \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{jl}(\lambda^t) &\leq \lambda_{ikjl}^{t+1} & \forall j, l, j < i, l \neq k \\ \lambda_{ikjl}^{t+1} &\leq \lambda_{ikjl}^t + \gamma_{ikjl}(\lambda^t) + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t) & \forall j, l, j < i, l \neq k.\end{aligned}$$

Therefore we have that the following inequalities imply constraints (26) and (27) for all $j, l, j > i, l \neq k$:

$$\mu_j'^{ik} + \nu_l'^{ik} \geq w_{ikjl} + \lambda_{ikjl}^t + \gamma_{ikjl}(\lambda^t) + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t)$$

and for all $j, l, j < i, l \neq k$

$$\mu_j'^{ik} + \nu_l'^{ik} \geq w_{ikjl} - \lambda_{jlik}^t + \gamma_{ikjl}(\lambda^t) + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t).$$

Constraints (26) and (27) are indeed implied, as, for all $j, l, j > i, l \neq k$,

$$\begin{aligned}\mu_j'^{ik} + \nu_l'^{ik} &= \mu_j^{ik} + \nu_l^{ik} + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t) \\ &\geq w_{ikjl} + \lambda_{ikjl}^t + \gamma_{ikjl}(\lambda^t) + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t)\end{aligned}$$

and for all $j, l, j < i, l \neq k$

$$\begin{aligned}\mu_j'^{ik} + \nu_l'^{ik} &= \mu_j^{ik} + \nu_l^{ik} + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t) \\ &\geq w_{ikjl} - \lambda_{jlik}^t + \gamma_{ikjl}(\lambda^t) + \left(\frac{1}{2(n_1 - 1)} + \frac{1}{2(n_2 - 1)} \right) \pi_{ik}(\lambda^t).\end{aligned}$$

Since $\mu_j^{ik}, \nu_l^{ik} \geq 0$ ($\forall j, l, j \neq i, l \neq k$) and by definition $\pi_{ik}(\lambda^t) \geq 0$, constraints (28) and (29) are satisfied as well. The objective value of (μ'^{ik}, ν'^{ik}) is given by

$$\sum_{\substack{j \\ j \neq i}} \mu_j'^{ik} + \sum_{\substack{l \\ l \neq k}} \nu_l'^{ik} = \sum_{\substack{j \\ j \neq i}} \mu_j^{ik} + \sum_{\substack{l \\ l \neq k}} \nu_l^{ik} + \pi_{ik}(\lambda^t) = v_{ik}(\lambda^t) + \pi_{ik}(\lambda^t).$$

Since (25) are minimization problems and there exist, for all i, k , feasible solutions with objective values $v_{ik}(\lambda^t) + \pi_{ik}(\lambda^t)$, we can conclude that the objective values of the solutions are bounded by this quantity. Hence, the lemma follows.

We use $\varphi = 0.5$, $\tau = 1$, and perform the dual descent method L successive times (see Algorithm 2).

Algorithm 2: DUALDESCENT(λ, L)

```

1  $\varphi \leftarrow 0.5$ ;  $[\text{LB}^*, \text{UB}^*] \leftarrow [Z_{\text{lb}}(\lambda), Z_{\text{LD}}(\lambda)]$ 
2 for  $n \leftarrow 1$  to  $L$  do
3   foreach  $i, k, j, l, i < j, k \neq l$  do
4      $\lambda_{ikjl} \leftarrow \lambda_{ikjl} + \varphi(\gamma_{ikjl} + \frac{\pi_{ik}(\lambda)}{2(n_1-1)} + \frac{\pi_{ik}(\lambda)}{2(n_2-1)}) - \varphi(\gamma_{jlik} + \frac{\pi_{jl}(\lambda)}{2(n_1-1)} + \frac{\pi_{jl}(\lambda)}{2(n_2-1)})$ 
5    $\text{LB}^* \leftarrow \max[\text{LB}^*, Z_{\text{lb}}(\lambda)]$ 
6    $\text{UB}^* \leftarrow Z_{\text{LD}}(\lambda)$ 
7 return  $[\text{LB}^*, \text{UB}^*]$ 

```

Overall method. Our overall method combines both the subgradient optimization and dual descent method. We do this performing the subgradient method until termination and then switching over to the dual descent method. This procedure is repeated K times (see Algorithm 3).

Algorithm 3: NATALIE(K, L, M, N)

```

1  $\lambda \leftarrow 0$ ;  $[\text{LB}^*, \text{UB}^*] \leftarrow [0, \infty]$ 
2 for  $k \leftarrow 1$  to  $K$  do
3    $[\text{LB}^*, \text{UB}^*] \leftarrow \text{SUBGRADIENTOPT}(\lambda, M, N) \cap [\text{LB}^*, \text{UB}^*]$ 
4    $[\text{LB}^*, \text{UB}^*] \leftarrow \text{DUALDESCENT}(\lambda, L) \cap [\text{LB}^*, \text{UB}^*]$ 
5 return  $[\text{LB}^*, \text{UB}^*]$ 

```

We implemented NATALIE in C++ using the LEMON graph library (<http://lemon.cs.elte.hu/>). The successive shortest path algorithm for maximum weight bipartite matching was implemented and contributed to LEMON. Special care was taken to deal with the inherent numerical instability of floating point numbers. Our implementation supports both the GraphML and GML graph formats. Rather than using one big alignment graph, we store and use a different alignment graph for every local problem (LD_λ^{ik}). This proved to be a huge improvement in running times, especially when the global alignment graph is sparse. NATALIE is publicly available at <http://planet-lisa.net>.

4 Experimental Evaluation

From the STRING database v8.3 [24], we obtained PPI networks for the following six species: *C. elegans* (cel), *S. cerevisiae* (sce), *D. melanogaster* (dme),

R. norvegicus (rno), *M. musculus* (mmu) and *H. sapiens* (hsa). We only considered interactions that were experimentally verified. Table 1 shows the sizes of the networks. We performed, using the BLOSUM62 matrix, an all-against-all global sequence alignment on the protein sequences of all $\binom{6}{2} = 15$ pairs of networks. We used affine gap penalties with a gap-open penalty of 2 and a gap-extension penalty of 10. The first experiment in Section 4.1 compares the raw performance of ISORANK, GRAAL and NATALIE in terms of objective value. In Section 4.2 we evaluate the biological relevance of the alignments produced by the three methods. All experiments were conducted on a compute cluster with 2.26 GHz processors with 24 GB of RAM.

species	nodes	annotated	interactions
cel (c)	5,948	4,694	23,496
sce (s)	6,018	5,703	131,701
dme (d)	7,433	6,006	26,829
rno (r)	8,002	6,786	32,527
mmu (m)	9,109	8,060	38,414
hsa (h)	11,512	9,328	67,858

Table 1: Characteristics of input networks considered in this study. The columns contain species identifier, number of nodes in the network, number of annotated nodes thereof, and number of interactions

4.1 Edge-Correctness

The objective function used for scoring alignments in GRAAL counts the number of mapped edges. Such an objective function is easily expressible in our framework using $s(a) = |\{(v, w) \in E_1 \mid (a(v), a(w)) \in E_2\}|$ and can also be modeled using the ISORANK scoring function. In order to compare performance of the methods across instances, we normalize the scores by dividing by $\min(|E_1|, |E_2|)$. This measure is called the edge-correctness by Kuchaiev et al. [17].

As mentioned in Section 3, our method benefits greatly from using a sparse alignment graph. To that end, we use the e-values obtained from the all-against-all sequence alignment to prohibit biologically unlikely matchings by only considering protein-pairs whose e-value is at most 100. Note that this only applies to NATALIE as both GRAAL and ISORANK consider the complete alignment graph. On each of the 15 instances, we ran GRAAL with 3 different random seeds and sampled the input parameter which balances the contribution of the graphlets with the node degrees uniformly within the allowed range of $[0, 1]$. As for ISORANK, when setting the parameter α —which controls to what extent topological similarity plays a role—to the desired value of 1, very poor results were obtained. Therefore we also sampled this parameter within its allowed range and re-evaluated the resulting alignments in terms of edge-correctness. NATALIE was

run with a time limit of 10 minutes and $K = 3$, $L = 100$, $M = 10$, $N = 20$. For both GRAAL and ISORANK only the highest-scoring results were considered.

Figure 2 shows the results. ISORANK was only able to compute alignments for three out of the 15 instances. On the other instances ISORANK crashed, which may be due to the large size of the input networks. For GRAAL no alignments concerning *sce* could be computed, which is due to the large number of edges in the network on which the graphlet enumeration procedure choked: in 12 hours only for 3% of the nodes the graphlet degree vector was computed. As for the last three instances, GRAAL crashed due to exceeding the memory limit inherent to 32-bit processes. Unfortunately no 64-bit executable was available. On the instances for which GRAAL could compute alignments, the performance—both in solution quality and running time—is very poor when compared to ISORANK and NATALIE. NATALIE outperforms ISORANK in both running time and solution quality.

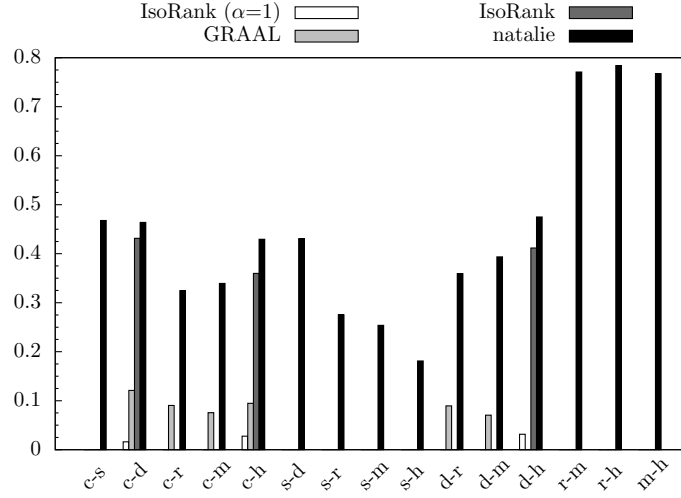
4.2 GO Similarity

In order to measure the biological relevance of the obtained network alignments, we make use of the Gene Ontology (GO) [3]. For every node in each of the six networks we obtained a set of GO annotations (see Table 1 for the exact numbers). Each annotation set was extended to a multiset by including all ancestral GO terms for every annotation in the original set. Subsequently we employed a similarity measure that compares a pair of aligned nodes based on their GO annotations and also takes into account the relative frequency of each annotation [11]. Since the similarity measure assigns a score between 0 and 1 to every aligned node pair, the highest similarity score one can get for any alignment is the minimum number of annotated nodes in either of the networks. Therefore we can normalize the similarity scores by this quantity. Unlike the previous experiment, this time we considered the bitscores of the pairwise global sequence alignments. Similarly to ISORANK parameter α , we introduced a parameter $\beta \in [0, 1]$ such that the sequence part of the score has weight $(1 - \beta)$ and the topology part has weight β . For both ISORANK and NATALIE we sampled the weight parameters uniformly in the range $[0, 1]$ and showed the best result in Figure 3. There we can see that both NATALIE and ISORANK identify functionally coherent alignments.

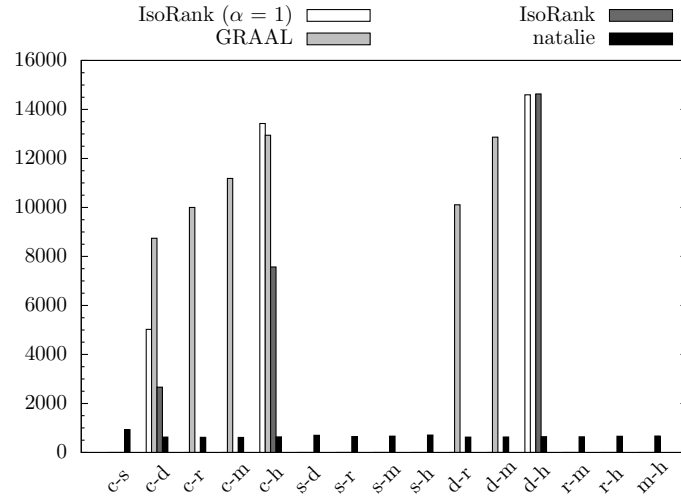
5 Conclusion

Inspired by results for the closely related quadratic assignment problem (QAP), we have presented new algorithmic ideas in order to make a Lagrangian relaxation approach for global network alignment practically useful and competitive. In particular, we have generalized a dual descent method for the QAP. We have found that combining this scheme with the traditional subgradient optimization method leads to fastest progress of upper and lower bounds.

Our implementation of the new method, NATALIE 2.0, works very well and fast when aligning biological networks, which we have shown in an extensive



(a) Edge correctness



(b) Running times in seconds

Fig. 2: Performance of the three different methods for the all-against-all species comparisons (15 alignment instances). Missing bars correspond to exceeded time/memory limits or software crashes.

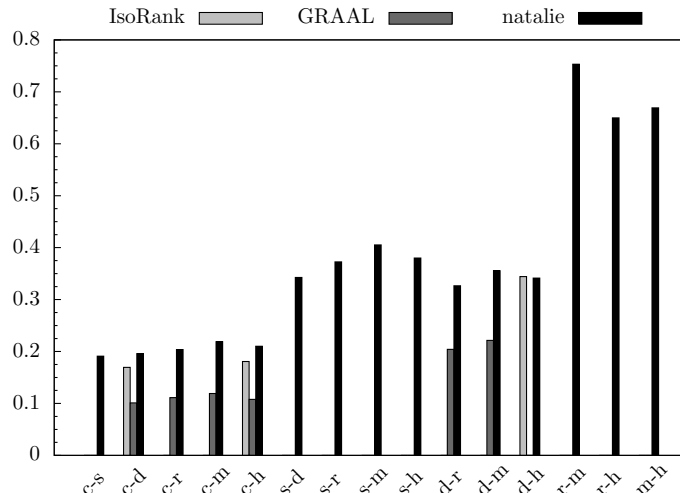


Fig. 3: Biological relevance of the alignments measured via GO similarity

study on the alignment of cross-species PPI networks. We have compared NATALIE 2.0 to those state-of-the-art methods whose scoring schemes can be expressed as special cases of the scoring scheme we propose. Currently, these methods are ISORANK and GRAAL. Our experiments show that the Lagrangian relaxation approach is a very powerful method and that it currently outperforms the competitors in terms of quality of the results and running time.

Currently, all methods, including ours, approach the global network alignment problem heuristically, that is, the computed alignments are not guaranteed to be optimal solutions of the problem. While the other approaches are intrinsically heuristic—both ISORANK and GRAAL, for instance, approximate the neighborhood of a node and then match it with a similar node—the inexactness in our methods has two causes that we plan to address in future work: On the one hand, there may still be a gap between upper and lower bound of the Lagrangian relaxation approach after the last iteration. We can use these bounds, however, in a branch-and-bound approach that will compute provably optimal solutions. On the other hand, we currently do not consider the complete bipartite alignment graph and may therefore miss the optimal alignment. Here, we will investigate preprocessing strategies, in the spirit of [25], to safely sparsify the input bipartite graph without violating optimality conditions.

The independence of the local problems (LD_{λ}^{ik}) allows for easy parallelization, which, when exploited would lead to an even faster method. Another improvement in running times might be achieved when considering more involved heuristics for computing the lower bound, such as local search. More functionally-coherent alignments can be obtained when considering a scoring function where node-to-node correspondences are not only scored via sequence similarity but also for instance via GO similarity. In certain cases, even negative weights for

topological interactions might be desired in which case one needs to reconsider the assumption of entries of matrix W being positive.

Acknowledgments. We thank SARA Computing and Networking Services (www.sara.nl) for their support in using the Lisa Compute Cluster. In addition, we are very grateful to Bernd Brandt for helping out with various bioinformatics issues and also to Samira Jaeger for providing code and advice on the GO similarity experiments.

References

1. W. P. Adams and T. Johnson. Improved linear programming-based lower bounds for the quadratic assignment problem. In *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, 1994.
2. U. Alon. Network motifs: theory and experimental approaches. *Nat Rev Genet*, 8(6):450–461, 2007.
3. M. Ashburner, C. A. Ball, J. A. Blake, et al. Gene ontology: tool for the unification of biology. *Nat Genet*, 25, 2000.
4. F. Ay, M. Kellis, and T. Kahveci. SubMAP: Aligning Metabolic Pathways with Subnetwork Mappings. *J Comput Biol*, 18(3):219–35, 2011.
5. A. Caprara, M. Fischetti, and P. Toth. A heuristic method for the set cover problem. *Oper Res*, 47:730–743, 1999.
6. J. Edmonds. Path, trees, and flowers. *Canadian J Math*, 17:449–467, 1965.
7. J. Edmonds and R. M. Karp. Theoretical improvements in algorithmic efficiency for network flow problems. *J ACM*, 19:248–264, 1972.
8. J. Flannick, A. Novak, B. S. Srinivasan, H. H. McAdams, and S. Batzoglou. Graemlin: general and robust alignment of multiple large interaction networks. *Genome Res*, 16(9):1169–1181, 2006.
9. M. Guignard. Lagrangean relaxation. *Top*, 11:151–200, 2003.
10. M. Held and R. M. Karp. The traveling-salesman problem and minimum spanning trees: Part II. *Math Program*, 1:6–25, 1971.
11. S. Jaeger, C. Sers, and U. Leser. Combining modularity, conservation, and interactions of proteins significantly increases precision and coverage of protein function prediction. *BMC Genomics*, 11(1):717, 2010.
12. M. Kanehisa, S. Goto, M. Hattori, et al. From genomics to chemical genomics: new developments in KEGG. *Nucleic Acids Res*, 34:D354–D357, 2006.
13. R. M. Karp. Reducibility among combinatorial problems. In R. E. Miller and J. W. Thatcher, editors, *Complexity of Computer Computations*, pages 85–103. Plenum Press, 1972.
14. B. P. Kelley, R. Sharan, R. M. Karp, et al. Conserved pathways within bacteria and yeast as revealed by global protein network alignment. *P Natl Acad Sci USA*, 100(20):11394–11399, 2003.
15. G. W. Klau. A new graph-based method for pairwise global network alignment. *BMC Bioinform*, 10 (Suppl 1):S59, 2009.
16. M. Koyutürk, Y. Kim, U. Topkara, et al. Pairwise alignment of protein interaction networks. *J Comput Biol*, 13(2):182–199, 2006.
17. O. Kuchaiev, T. Milenkovic, V. Memisevic, W. Hayes, and N. Przulj. Topological network alignment uncovers biological function and phylogeny. *J R Soc Interface*, 7(50):1341–54, 2010.

18. H. W. Kuhn. The Hungarian method for the assignment problem. *Nav Res Logist Q*, 2(1-2):83–97, 1955.
19. E. L. Lawler. The quadratic assignment problem. *Manage Sci*, 9(4):586–599, 1963.
20. J. Munkres. Algorithms for the assignment and transportation problems. *SIAM J Appl Math*, 5:32–38, 1957.
21. R. Sharan and T. Ideker. Modeling cellular machinery through biological network comparison. *Nat Biotechnol*, 24(4):427–433, 2006.
22. R. Sharan, T. Ideker, B. Kelley, R. Shamir, and R. M. Karp. Identification of protein complexes by comparative analysis of yeast and bacterial protein interaction data. *J Comput Biol*, 12(6):835–846, 2005.
23. R. Singh, J. Xu, and B. Berger. Global alignment of multiple protein interaction networks with application to functional orthology detection. *P Natl Acad Sci USA*, 105(35):12763–12768, 2008.
24. D. Szklarczyk, A. Franceschini, M. Kuhn, et al. The STRING database in 2011: functional interaction networks of proteins, globally integrated and scored. *Nucleic Acids Res*, 39:561–568, 2010.
25. I. Wohlers, R. Andonov, and G. W. Klau. Algorithm engineering for optimal alignment of protein structure distance matrices. *Optim Lett*, 2011.